

Development of a Multilingual Information Retrieval and Check System Based on Database Semantics*

*Kiyong Lee (KOREA U), Suk-Jin Chang (SEOUL N. U),
Yun-Pyo Hong (DANKUK U), Key-Sun Choi (KAIST),
Minhaeng Lee (YONSEI U), Jae Sung Lee (ETRI, KOREA),
Jungha Hong (KOREA U), Juho Lee (KAIST),
Junsik Hong (YONSEI U)*

1 Introduction

Technical documents for multilateral agreements or international business transactions are normally produced in a bilingual or multilingual form. Being mostly of legal nature, these documents require especially accurate and speedy translations by expert translators. In order to aid these experts, automatic ways of checking translation results (such as a spelling checker) would be highly desirable.

This paper describes the MIRAC system for Multilingual Information Retrieval And Checking. It is designed to find translation errors in multilingual documents, and to evaluate the overall results of translation. Unlike a machine translation or a translation memory system [Volk98, Webb98], the primary function of the MIRAC system is to evaluate previously translated and aligned documents in source and target languages, while dynamically building a database that consists of aligned multilingual texts carrying semantically equivalent content.

MIRAC consists of two components: one is a lexical evaluation module and the other is a semantic evaluation module, based on Hausser's Database Semantics [Hausser99].¹ Instead of aiming at the metric evaluation of machine translation systems, MIRAC directly evaluates translated documents by checking first the consistency of use of lexical terms and then semantic equivalences between source and target documents.

The paper is organized as follows: a brief introduction of Termight, a workbench for technical translators, in section 2; an introductory overview of the MIRAC system, with a description of each of its parts, in section 3; a report on its implementation and experiment in section 4; and concluding remarks in the final section 5.

2 Related Works: Termight

Dagan and Church's Termight [Dagan/Church94] is a workbench for technical translators. It mainly checks the correctness of translated technical terminology. The process is semi-automatic, for it requires a manual listing of technical terms in an original text before their corresponding translations are automatically searched from the translated text.

For listing technical terms, Termight analyzes a document for part-of-speech tagging and finds compound nouns. Out of these compound nouns, technical terms are identified and edited appropriately under a suitable environment provided by the system. Termight's alignment program then automatically locates their corresponding translations, while correct translations are selected manually to build a translation glossary. Here, the workbench helps to find correct translation pairs.

3 Overall Structure of the MIRAC System

The overall architecture of the MIRAC system is shown below.

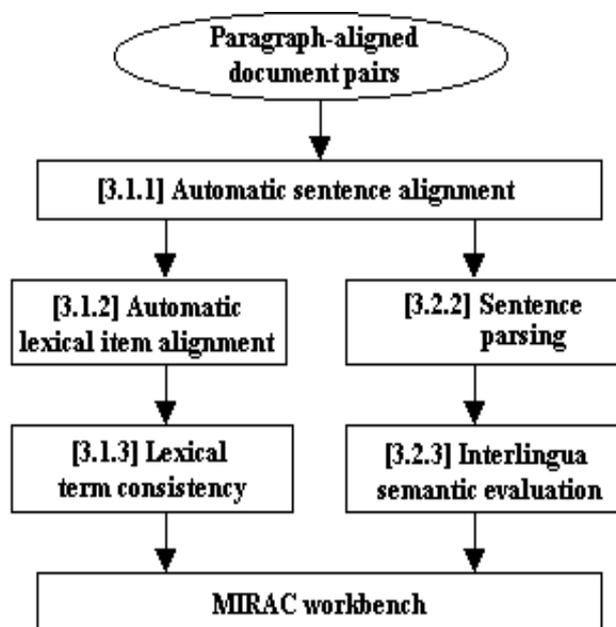


Figure 1: Structure of the MIRAC system.

The MIRAC system deals with multilingual documents, specifically comparing a pair of documents in a source and a target language. As its input, MIRAC takes in paragraph-aligned document pairs in these two languages [Klee/Park97]. Then these pairs of documents are aligned at both the sentential and lexical levels by an automatic alignment program. The use consistency of technical or key terms is also automatically checked by another module. These processes are carried out by a statistical method as pre-processing steps for the evaluation of correct translation [Collier/Ono/ Hira98].

Translations are evaluated in two steps. The first step evaluates the lexical correspondence between pairs of the aligned documents, displaying the results of evaluation in the alignment workbench. The system checks the correctness and consistency of the use of translated terms in the target language. The second step checks the semantic correspondence by a statistical method between the corresponding pairs of terms, phrases, and sentences, again displaying the results of evaluation on the alignment workbench.

3.1 Lexical Evaluation

Lexical items in each pair of aligned sentences are all aligned automatically by a statistical method [Jslee/Kang/Jhlee/Le/Choi97]. The module for lexical evaluation then analyzes them and displays the two lists of original and translated terms on its workbench. The correctness of translation should, however, be checked manually.

3.1.1 Automatic Sentence Alignment

For lexical item alignment, sentences must be aligned first. For this, the Gale/Church method [Gale/Church91] is used to measure the length of each sentence for statistical calculation. This method, however, needs to be improved by providing ways of using information from dictionaries and also from the feedback of evaluation processes.

For our experiment, the introductory chapter of Negroponte's (1995) *Being Digital*, both in its original English and in its Korean translation, was analyzed. Each sentence and paragraph in the chapter was marked for the experiment. Both the English and the Korean versions were found to contain 20 paragraphs each. The English version contains 81 sentences, but the Korean version contains 86.

The accuracy of alignment is 97.47%. In this short experiment, the statistical method was found to be very fast and efficient but produced a far less satisfactory result on alignment accuracy. For its improvement, the additional use of dictionary information should be helpful [Collier/Ono/Hira98].

3.1.2 Automatic Lexical Item Alignment

Lexical items are statistically aligned on the basis of co-occurrence information [Hull98, Jhlee99]. The overall process is shown in Figure 2.

The intermediate steps are as follows:

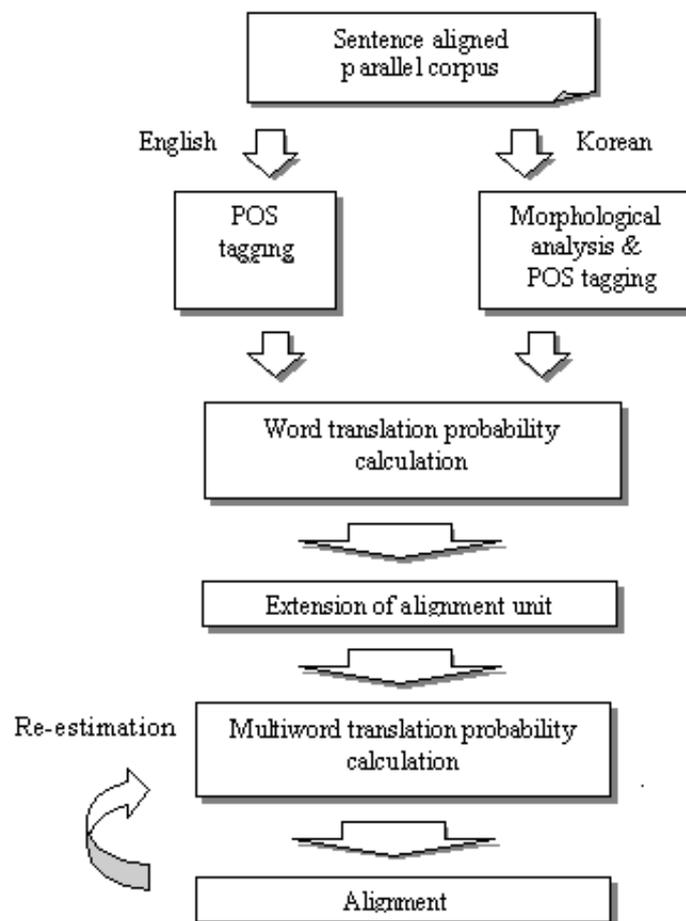


Figure 2: Overall flow of automatic alignment.

(i) For easy alignment, content words are extracted from each language document. In English, content words are nouns, verbs, or adjectives. In Korean, however, only nouns and noun-derived verbs or adjectives are treated as content words because pure verbs and adjectives are rarely used as technical terms.

(ii) The probability of word translation is calculated on the basis of information on bilingual co-occurrence. The basic assumption is that the word translation probability is higher if a word and its translation occur more frequently in aligned sentence pairs. The calculation of translation probability or similarity between two words is usually based on their respective meaning and information as well as their Dice coefficients that provide co-occurrence information.

In this experiment, Dice coefficients are used to calculate word translation probability. The translation probability $C_p(E_i, K_j)$ for an English word E_i and its corresponding Korean word K_j , for instance, is defined as follows:

$$C_p(E_i, K_j) = \frac{2C(E_i, K_j)}{C(E_i) + C(K_j)}$$

$C(E_i)$: the number of segments which contain E_i

$C(E_i, K_j)$: the number of pair segments which contain E_i and K_j in each segments

(iii) This method of calculating word translation probabilities can be extended to the calculation of multiword translation probabilities by including neighboring content words among their alignment units. It is assumed here that sequences of neighboring words can be formed into multi-content words. For an easy implementation, the following four cases are considered here: 1:1, 1:2, 2:1, and 2:2 types of word correspondence.

Each of the 1:2 and 2:1 cases can be extended only if all of the following conditions are satisfied:

$$C_p(E_i, K_j K_{j+1}) \geq \frac{C_p(E_i, K_j) + C_p(E_i, K_{j+1})}{2}$$

$$C_p(E_i E_{i+1}, K_j) \geq \frac{C_p(E_i, K_j) + C_p(E_{i+1}, K_j)}{2}$$

or

$$C_p(E_i, K_j K_{j+1}) \geq \max\{C_p(E_i, K_j), C_p(E_i, K_{j+1})\}$$

$$C_p(E_i E_{i+1}, K_j) \geq \max\{C_p(E_i, K_j), C_p(E_{i+1}, K_j)\}$$

The 2:2 case is a little more complicated. It can be extended only if the following condition is satisfied:

$$\begin{aligned}
 C_p(E_i E_{i+1}, K_j K_{j+1}) &\geq \frac{C_p(E_i, K_j) + C_p(E_{i+1}, K_{j+1})}{2} \text{ and} \\
 &\geq \frac{C_p(E_i, K_{j+1}) + C_p(E_{i+1}, K_j)}{2} \text{ and} \\
 &\geq \frac{C_p(E_i E_{i+1}, K_j) + C_p(E_i E_{i+1}, K_{j+1})}{2} \text{ and} \\
 &\geq \frac{C_p(E_i, K_j K_{j+1}) + C_p(E_{i+1}, K_j K_{j+1})}{2} \\
 &\text{or} \\
 C_p(E_i E_{i+1}, K_j K_{j+1}) &\geq \max\{C_p(E_i, K_j), C_p(E_i, K_{j+1}), C_p(E_{i+1}, K_j), C_p(E_{i+1}, K_{j+1}), \\
 &C_p(E_i E_{i+1}, K_j), C_p(E_i E_{i+1}, K_{j+1}), C_p(E_i, K_j K_{j+1}), C_p(E_{i+1}, K_j K_{j+1})\}
 \end{aligned}$$

A distance limit should be imposed to exclude meaningless multi-words that form parts of a content word but are separated by too great a distance.

3.1.3 Lexical Term Consistency Checking

For the accuracy and quality of translation, the consistent use of terms should be checked, especially in technical documents. This is especially so in the case of technical terms and proper nouns; otherwise only confusion will arise.

For example, the term “computer” is usually translated to “khem.phyu.the”, but it can be translated to “khom.phyu.the”, “cen.ca.kyey.san.ki”, “cen.san.ki”, and so on.² Someone familiar with the concept of a computer may think they are all the same, but others may not, for “cen.ca.kyey.san.ki” normally refers to a calculator. Another example is the name of a university in Chonju, Korea. It has a unique Korean name that has been translated or, more accurately speaking, romanized into Jeonbug, Chonpuk, or Chonbuk National University, causing great confusion.

In order to evaluate lexical consistency, we need to list lexical items and their corresponding translations. The MIRAC workbench first extracts them from a pair of documents, producing an aligned lexical list. It then analyzes the list to examine the consistency of their use. The results of these analyses can be used to build a translation dictionary for further use in checking the accuracy of translation.

3.2 Semantic Evaluation

The semantic evaluation module of MIRAC requires the parsing of each pair of aligned sentences from both source and target languages. Being implemented within the Malaga system, it analyzes each sentence left-associatively, yielding its result in an attribute-value matrix (AVM) form. These AVMs contain semantic information that constitutes an interlingua (IL), thus allowing the bi-directional translation of one language to another. The semantic evaluation of each pair of source and target sentences is then carried on by comparing their semantic information in the interlingua.

3.2.1 A Theoretical Basis

For semantic evaluation, MIRAC adopts Hausser's Database Semantics [Hausser99].³ It allows the representation of propositional content and other related semantic information in an abstract interlingua, thus making it possible to evaluate both the consistency of the TL-formulations and their adequacy vis à vis the propositional content.

The technical basis of this evaluation is the transition counters characteristic of database semantics. In current systems, transition counters indicate which navigation through the propositional content is the most recent and which navigations are the most frequent. The purpose of the counters is to ensure that the autonomous navigation underlying conceptualization in language production proceeds without splits and loops.

It is conceivable, however, to employ counters in other applications as well. For example, in order to model the learning of new fashionable formulations, additional counters would be implemented at the level of natural language.

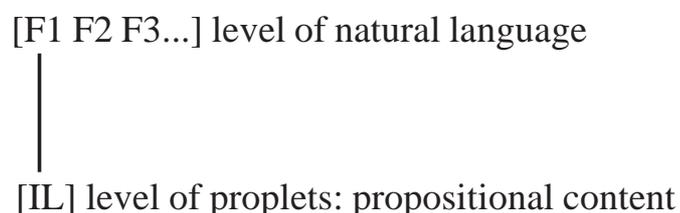


Figure 3: Matching of formulations in NL and proplet levels.

Here alternative formulations, F1, F2, F3, etc., for the same propositional content have values for their frequency in interpretation and production. Based on these frequency values, the speaker could choose a common or a special formulation depending on the utterance situation.

In IL-based MT, the above schema is extended as follows.

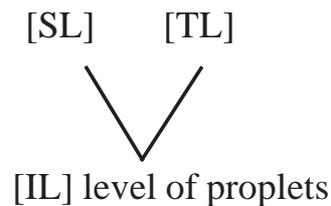


Figure 4: Convergence of SL and TL formulations at the level of proplets.

This schema suggests another possible use of counters: they mark not only alternative formulations of the SL and the TL for frequency relative to corresponding IL propositions, but also their correlation to each other.

The characteristic technical environment of database semantics is especially suited for an efficient implementation of counters. Furthermore, database semantics is special because it treats (i) the IL formally as an unordered set of AVMs and (ii) the interpretation and production procedures alike on the common basis of a time-linear navigation. These structural properties are ideally suited for storing the information specific to translation memory or database.

3.2.2 Sentence Parsing

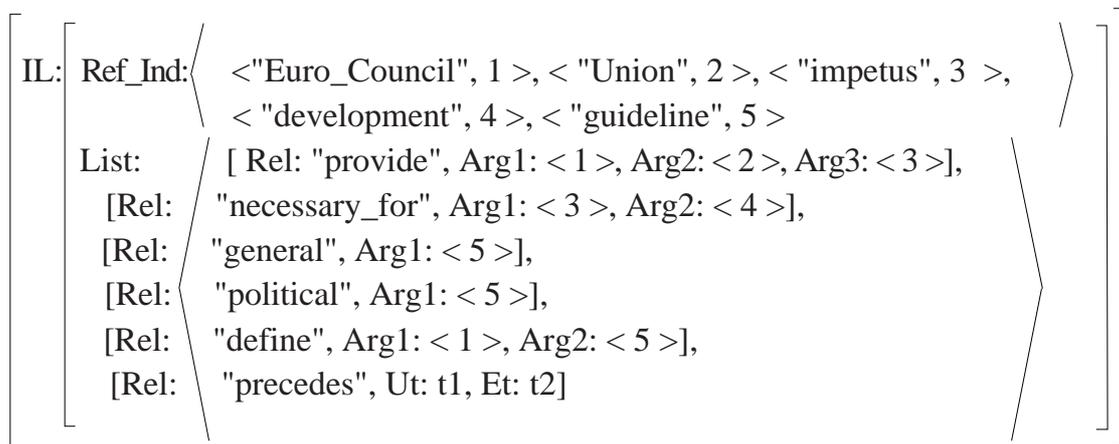
Although both source and target languages have their own distinct systems, these systems have the same structure with the same theoretical basis, namely Hausser's Left-Associative Grammar [Hausser99], and are all implemented in the same programming language, Malaga. As for Korean, for instance, Lee [Klee99a] implemented its morphological analyzer Komor and Hong/Lee [Hong/ Klee99] its syntactic parser.

3.2.3 Interlingua Semantic Evaluation

In order to allow a bidirectional translation from one language to another and also to evaluate its correspondence, sentences are mapped into the interlingua format of sets of proplets. These are created in the process of morphological and syntactic parsing.⁴

Each pair of parsed sentences in the source and target languages with their semantic content is now checked with the evaluation module for semantic equivalence or identity. The following example shows how the semantic content of sentence (1) is represented in Interlingua.

- (1) The European Council shall provide the Union with the necessary impetus for its development and shall define the general political guidelines.



The attribute *IL* takes as value two complex features, *Ref_Ind* and *List*. The first feature consists of an attribute *Ref-Ind* and its value that lists all of the key terms occurring in sentence (1). The second feature, on the other hand, simply consists of a list of proplets, each representing basic propositional content conveyed by the sentence. In each proplet, a relation *Rel* takes more than one argument *Arg*, while each *Arg* is related through an index number to a key term listed in *Ref_Ind*. The last proplet states that the utterance time *t1* precedes the event time *t2*, thus referring to an event occurring in the future.

Assuming that we have obtained a similar, if not the same, matrix representation for a Korean translation of sentence (1), the evaluation module checks and gives an evaluation point for each of the following items:⁵

(2) Evaluation Items and Scores

<i>items</i>	<i>scores</i>	
·Proposition-Relation	40 (35)	
·Reference-Indices	30 (25)	
·Modification	20 (10)	
·Tense	10 (10)	
	100 (80)	→ very good

4 Implementation and Experiment

4.1 Experiment of Lexical Evaluation

The English-Korean parallel corpus, consisting of a 750-page volume on Uruguay Round multilateral agreements, was used for our experiments. The corpus is aligned in segment units in the preprocessing step. Each segment is mostly composed of one single sentence. Some statistical facts about the parallel corpus are given below:

items	English	Korean
segments	4,968	4,968
words (phrases)	139,265	79,290
average length of segments	28.03	15.96
content words	65,844	65,653
unique content words	2,681	3,847

Table 1: Statistics for the parallel corpus.

This table shows that the number of Korean words or word groups occurring in the corpus is smaller than that of English words because compound words are used more frequently in Korean. The average number of words occurring in each

of the English sentences is 28.03, while the average number of words occurring in each of the Korean sentences is 15.96.

The experiment was performed in several steps. Each language document from the parallel corpus is POS-tagged by a language tagger. The tagging process filters the content words. For tagging English, [Brill94]'s method was used. For Korean, an English HMM (Hidden Markov Model) tagger was modified to process Korean sentences [Shin/Han/ Park/Choi95].

The translation probability of content words is calculated on the basis of bilingual co-occurrence information. By extending it to their neighboring words, the translation probability of multi-words is then calculated. This probability is used to align the multi-words and then is recalculated by counting the aligned multi-words.

This process is normally repeated seven times. In order to find meaningful multi-words, positional information is also used. The following graphs (Figure 5) show the change of the number of extracted unique translation pairs and the percentage of each type of correspondence at the last alignment. We can see that they both show a similar trend: the number of translation pairs decreases rapidly at the first re-estimation but decreases very slowly from the second re-estimation and finally remains constant after the fourth re-estimation.

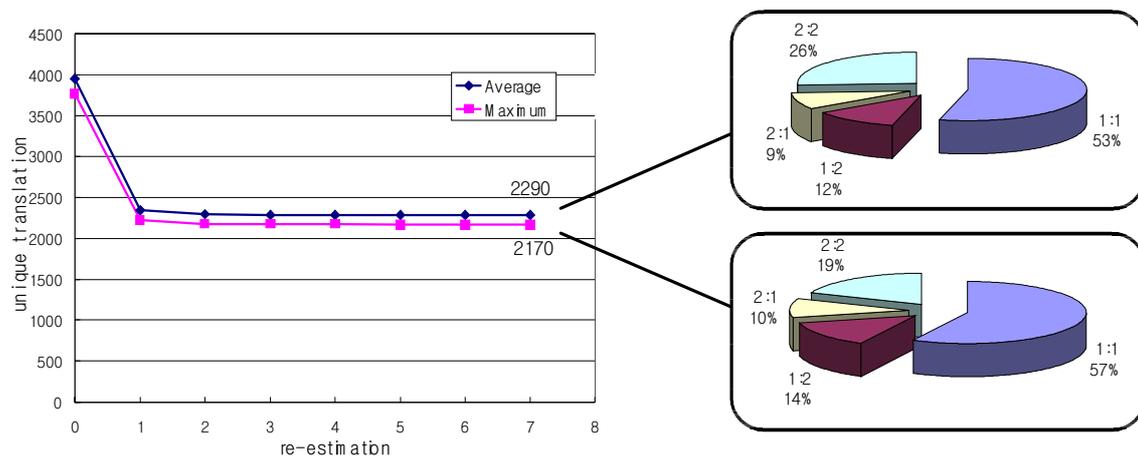


Figure 5: Change in the number of extracted unique translation pairs.

After the fourth re-estimation, there is also no change in translation pairs, while there is a small change in the translation probability. We assume that both of the results converge at the seventh iteration. After the seventh re-estimation, the ratio of the types of corresponding pairs in each case is similar. The number of correspondences of type 1:1 is the largest, type 2:2 the second largest, type 1:2 the next, and type 2:1 the last.

For the measurement of alignment accuracy, two methods were used: the calculation average method for distance limits and the maximum value method. In this experiment, the accuracy was calculated for 100 randomly selected translation pairs by using both methods. They both produced similar test results. The accuracy of alignment at each re-estimation is shown in Figure 6.

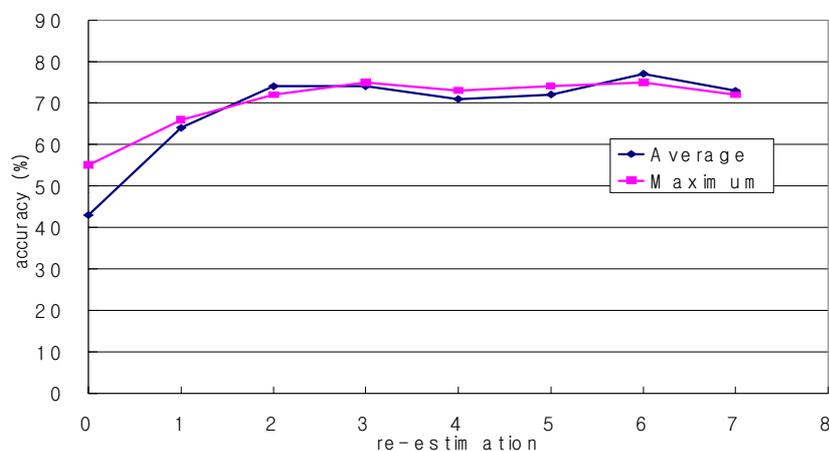


Figure 6: Accuracy of alignment.

The accuracy increases at the beginning of iteration, but stops changing after a certain point. The two curves in Figure 4 show almost the same results. The calculation average method, for instance, produced 2,290 unique translation pairs at the convergence points (from the fifth iteration to the seventh) with an average accuracy of 74%.

Some examples of the extracted translation pairs are shown in Table 2.

	English	Korean
(3)	convention	hyep.yak
(4)	countermeasure	tay.ung co.chi
(5)	working party	cak.ep.pan
(6)	result (of) negotiation(s)	hyep.sang keyl.kwa
(7)	result (of) negotiation(s)	hyep.sang(.uy) keyl.kwa

Table 2: Examples of extracted translation pairs.

The experiment was performed on the Uruguay Round (UR) documents on the world economy and diplomatic affairs. Since these documents deal with problems in a very specialized domain, no translation dictionary of general use provides appropriate translation words. A case in point is the pair “convention” : “hyep.yak” shown in (3). The English word “convention” generally means “cip.hoy” (meeting), “kwan.lyey” (traditional case), “sa.hoy.cek. kwan.swup” (social custom) in Korean. In these documents, “convention” is aligned to the diplomatic term “hyep.yak”.

Here, various types of alignment were found. The alignment “convention” : “hyep.yak” in (3) is of type 1:1, while the alignment “countermeasure” : “tay.ung co.chi” in (4) is of type 1:2. The alignment “working party” : “cak.ep.pan” in (5), on the other hand, is of type 2:1. One word or phrase may have two different alignments, too: for example, the phrase “result (of) negotiation(s)” is aligned to “hyep.sang kyel.kwa” in (6) or to “hyep.sang.uy kyel.kwa” in (7). But these two are the same if only content parts are taken into account, for the particle “uy” in Korean is a function word meaning “of”. This shows that our method of extending multi-words is effective for finding content multi-words in various lexical forms.⁶

Figure 7 (following page) shows a screen shot of the MIRAC workbench. If a user selects a word in a source language, its translation candidates and their translation probabilities are shown in the workbench. When the user selects one of the translation candidates, some translation examples appear in the main screen. The screen shows the word “committee” being translated to “wi.wen.hoy”, displaying its translation candidates and a list of translation examples.

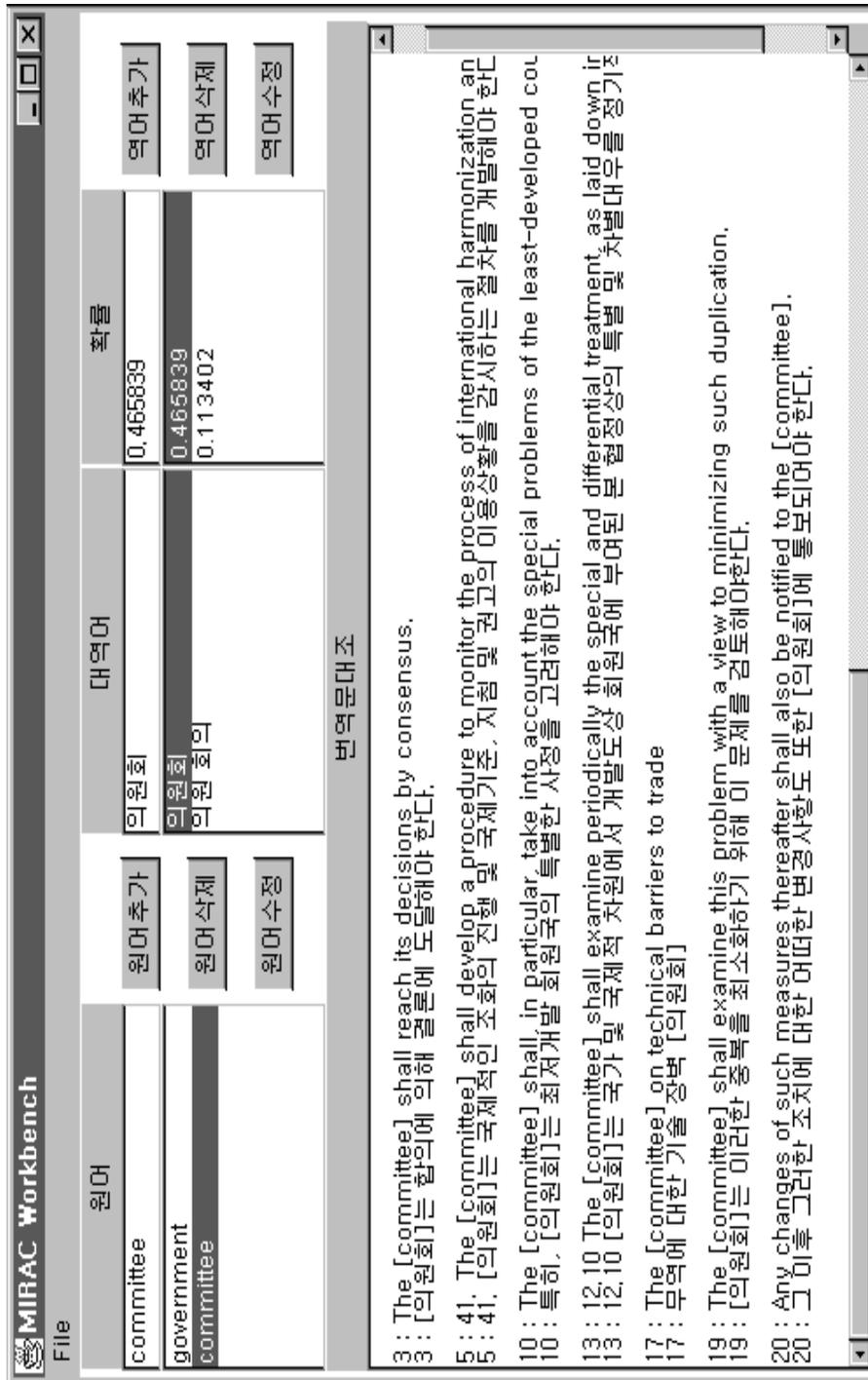


Figure 7: Screenshot of the MIRAC system.

4.2 Experiment of Semantic Evaluation

Here, we randomly selected five English sentences from the UR documents and translated them into two different languages, Korean and German. Then these translations were compared to check how their semantic content was preserved in each of the translations. Table 3 shows the results of the comparison.

	TL : Korean		TL : German	
	Scores	Violations	Scores	Violations
Sentence 1	75	Prop-Rel, Ref-Ind	100	
Sentence 2	85	Prop-Rel	80	Ref-Ind, Mod
Sentence 3	100		100	
Sentence 4	75	Prop_Rel, Ref-Ind	70	Ref-Ind, Mod
Sentence 5	90	Ref-Ind	100	
Average	85		90	

Table 3: Results of semantic evaluation (Korean and German).

The average scores in Table 3 show that the German translation scored higher than the Korean translation. While the Korean translation often failed to capture Prop(osition)-Rel(ation) and Ref(erence)-Ind(ices), the German translation failed to capture Ref-Ind and Mod(ification) in sentences 2 and 4. These results indicated that the translation between typologically similar SL and TL, like English and German, was easier than the translation between typologically dissimilar languages like English and Korean. The latter case even failed to capture such basic relations like Prop-relations.

5 Concluding Remarks

Machine translation is a formidable task. The task of evaluating the results of translation, however, is more tractable. The MIRAC system shows such a possibility by demonstrating how multilingual texts can be systematically aligned for checking the consistency of the use of lexical terms as well as the semantic equivalences between source and target languages.

The MIRAC system evaluates the quality of previously translated documents aligned in source and target languages, while continuously updating a database consisting of such aligned multilingual texts. It thus closely resembles a translation memory system. Nevertheless, its main function is to systematically evaluate the accuracy of translations at both the lexical and the propositional level. An evaluation tool like the MIRAC system is not only useful, but also necessary for building an adequate translation memory or storage system as well as an efficiently running machine translation system. When combined into one coherent system, these three systems of evaluation, memory, and translation can become a constantly or dynamically upgrading integrated machine translation system.

Especially when source and target languages differ from each other structurally, the evaluation of semantic equivalence plays an important role. This, for instance, should be the case, when a non-western language like Korean or Chinese is translated into English or vice versa. Being based on Hausser's Database Semantics [Hausser99], the MIRAC system can adequately represent the semantic content of sentences in both source and target languages in terms of abstract proplets and check their semantic equivalence. Contents, stored in the MIRAC system, can be recycled to evaluate both the consistency of a TL-formulation and its adequacy relative to the propositional content.

Since it is based on Database Semantics, the MIRAC system can also be implemented to be part of a machine translation system. For it can reproduce acceptable sentences in a target language by navigating through a word bank or an arrayed field of proplets that has been built of a source language and then by selecting appropriate sequences of words or proplets. In further research, Database Semantics may thus be extended into an approach to machine translation where translation memory serves not only as an aid to human translation, but as an important component of automatic translation.

References

- [Beutel97] Beutel, Bjoern (1997): 'Malaga 4.3', Abteilung Computerlinguistik, Universität Erlangen-Nürnberg, <http://www.linguistik.uni-erlangen.de/Malaga.de.html>.
- [Brill94] Brill, Eric (1994): 'Some Advances in Transformation-Based Parts of Speech Tagging', *Proceedings of the 12th National Conference of Artificial Intelligence*, 722–727.

- [**Chang98**] Chang, Suk-Jin (1998): ‘Translation: Mapping and Evaluation’, [written in Korean], *Language and Information 2.1*, 1–41.
- [**Collier/Ono/Hira98**] Collier, Nigel, Kenji Ono, and Hideki Hirakawa (1998): ‘An Experiment in Hybrid Dictionary and Statistical Sentence Alignment’, *Proceedings of the 17th International Conference on Computational Linguistics*, 268–274.
- [**Cop/Flick/Mal/Riehe/Sag96**] Copestake, Ann, Dan Flickinger, Robert Malouf, Susanne Riehemann, and Ivan A. Sag (1996): ‘Translation using Minimal Recursion Semantics’, *Proceedings of the Sixth International Conference on Theoretical and Methodological Issues in Machine Translation*, Leuven.
- [**Dagan/Church94**] Dagan, Ido and Ken Church (1994): ‘Termight: Identifying and Translation Technical Terminology’, *Proceedings of the 4th Conference on Applied Natural Language Processing*, 34–40.
- [**Gale/Church91**] Gale, William A. and Kenneth W. Church (1991): ‘A Program for Aligning Sentences in Bilingual Corpora’, *Proceedings of the 29th Annual Meeting of the Association for Computational Linguistics*, 177–184.
- [**Hausser99**] Hausser, Roland (1999): *Foundations of Computational Linguistics: Man-Machine Communication in Natural Language*, Berlin: Springer-Verlag.
- [**Hong/Klee99**] Hong, Jungha and Kiyong Lee (1999): ‘Processing Korean Relative Adnominal Clauses’, [written in Korean], *Proceedings of the Eleventh Conference on Korean Characters and Korean Information Processing*, 265–271.
- [**Hull98**] Hull, David A. (1998): ‘A Practical Approach to Terminology Alignment’, *Proceedings of the First Workshop on Computational Terminology*, 1–7.
- [**Jhlee99**] Lee, Juho (1999): *Extraction of English-Korean Compound Noun Translation through Automatic Method*, [written in Korean], Master’s thesis, KAIST.
- [**Jslee/Kang/Jhlee/Le/Choi97**] Lee, Jae Sung, Jung-Gu Kang, Juho Lee, Hung Le, Key-Sun Choi (1997): ‘Design and Implementation of Alignment Workbench’, [written in Korean], *Proceedings of the Ninth Conference on Korean Characters and Korean Information Processing*, 430–435.

- [**Klee/Park97**] Lee, Kiyong and Chongwon Park (1997): ‘On Designing a Term Match Checking System for Bilingual Documents’, [written in Korean], *Proceedings of the 1997 Spring Conference on Cognitive Science*, 220–229.
- [**Klee99a**] Lee, Kiyong (1999a): *Computational Morphology*, [written in Korean], Seoul: Korea University Press.
- [**Klee99b**] Lee, Kiyong (1999b): ‘A Basis of Database Semantics: from Feature Structures to Tables’, [written in Korean], *Proceedings of the Eleventh Conference on Korean Characters and Korean Information Processing*, 297–303.
- [**Lee/Jee/Chung98**] Lee, Minhaeng, Kwangsin Jee, So W. Chung (1998): ‘A Research on Test Suites for Translation Systems’, [written in Korean], *Language and Information 2.2*, 185–220.
- [**Shin/Han/Park/Choi95**] Shin, Jung H., Young S. Han, Young C. Park, and Key S. Choi (1995): ‘A HMM Part-of-Speech Tagger for Korean with Word-Phrasal Relations’, *Proceedings of Recent Advances in Natural Language Processing*, 439–449.
- [**Volk98**] Volk, Martin (1998): ‘The Automatic Translation of Idioms, Machine Translation vs. Translation Memory Systems’, in: Nico Weber(ed.), *Machine Translation: Theory, Applications, and Evaluation. An Assessment of the State of Art*, St. Augustin: Gardez Verlag.
- [**Webb98**] Webb, Lynn E. (1998): ‘Advantages and Disadvantages of Translation Memory: A Cost/Benefit Analysis’, MA Thesis, Monterey Institute of International Studies.

ENDNOTES

- * This paper was supported by a NON-DIRECTED RESEARCH FUND, Korea Research Foundation, 1996. This is the final report on a three-year (1996-1999) joint collaborative research project (PI: Prof. Kiyong Lee, Korea University) with overseas consultant Prof. Roland R. Hausser and his research staff at the University of Erlangen-Nürnberg, Germany. We are grateful to these organizations for their financial and technical support. We are also grateful to Chongwon Park, Dr. Kunsik Lee, Dr. Kung-Un Choi, Dr. Si-Jong Ryu, and Koaunghi Un, who participated in the project as research associates in the past two or three years. We also owe our gratitude to the anonymous referees for their constructive comments, Prof. Uta Seewald-Heeg and

Rita Nübel for their critical review and editorial assistance, and finally Mr. Gary Rector for his professional styling and proofreading.

- ¹ This is an extension of Hausser's Left-Associative Grammar [Hausser99], implemented in Beutel's C-like programming language called Malaga [Beutel97], which accommodates LAG with attributes.
- ² In this paper, Hangul is romanized using the Yale system.
- ³ Lee [Klee99c] proposed a slightly different version of Database Semantics by adopting an object-oriented relational model, for it can easily convert AVMs for natural language into table forms and allows the use of SQL for developing a natural language query system.
- ⁴ Here we try to adopt Copestake et al.'s representation schema [Cop/Flick/Mal/Riehe/Sag96], which is introduced in their Minimal Recursion Semantics.
- ⁵ A more detailed scheme of evaluation for the MIRAC system is presented in [Chang98] and [Lee/Jee/Chung98].
- ⁶ Just as inflectional endings or prepositions rarely carry any content in English, nominal particles like "uy" carry practically no content in an agglutinative language like Korean or Japanese.